

# 修士論文の和文要旨

大学院情報システム学 研究科		博士前期課程	情報システム設計学	専攻
氏名	矢部 走		学籍番号	0650035
論文題目	文書構造化を用いた Web ページからの RSS 自動生成システム			
要 旨				
<p>多くの Web ページから必要な情報を効率的に取得するための手段として RSS が注目を集めている。RSS は Web サイトの各ページのタイトル，URL，概要，更新時刻などを構造化して記述する XML ベースのフォーマットで，Web ページの更新情報を配信するのに適しており，多くのサイトで提供され始めている。しかし，現状では RSS を配信していない Web ページも多く RSS の利用の弊害となっている。その問題を解決するため，Web ページから RSS を自動生成するサービスもいくつか存在する。それらのサービスで用いられているアルゴリズムは日付抽出とリンク抽出の 2 つに大別できるが，どちらも対応するページの種類が限られており，RSS を生成することが出来ない場合も多い。特に，テキスト主体の Web ページにおける RSS の生成は難しいのが現状である。</p> <p>そこで本論文では，テキスト主体の Web ページからも RSS を生成すべく，既存の手法である文書構造化，そして日付表現を用いたタイトル抽出を用い，さらにそれに独自の要素として HTML 文書をグループ構造化し，そのグループ構造と DP マッチングを用いてタイトルの抽出を行う方法を組み合わせた RSS 生成手法を提案し，従来 RSS 情報の抽出が困難であった Web ページからの RSS 自動生成を試みる。</p> <p>実験では 18 種類，153 のページを対象に既存手法と提案手法それぞれによる RSS 生成を行い，精度約 65 % の既存手法に対し，提案手法は約 85 % の精度で RSS を生成できることを示した。さらに提案手法は，既存の手法では RSS 生成が難しかったテキスト系の Web ページからの RSS 生成に対応し，他の種類の Web ページよりも精度が高いという結果を残した。</p>				